



PCIe® Re-timer in Data Centers Platforms

Alex Umansky
Chief Architect, PCIe Solutions
IT PL, Huawei

Disclaimer



Presentation Disclaimer: All opinions, judgments, recommendations, etc. that are presented herein are the opinions of the presenter of the material and do not necessarily reflect the opinions of the PCI-SIG®.

Agenda



- **Introduction**
- **System Integration Challenges**
- **Re-timer Affect on Latency and Performance**
- **Conclusions**

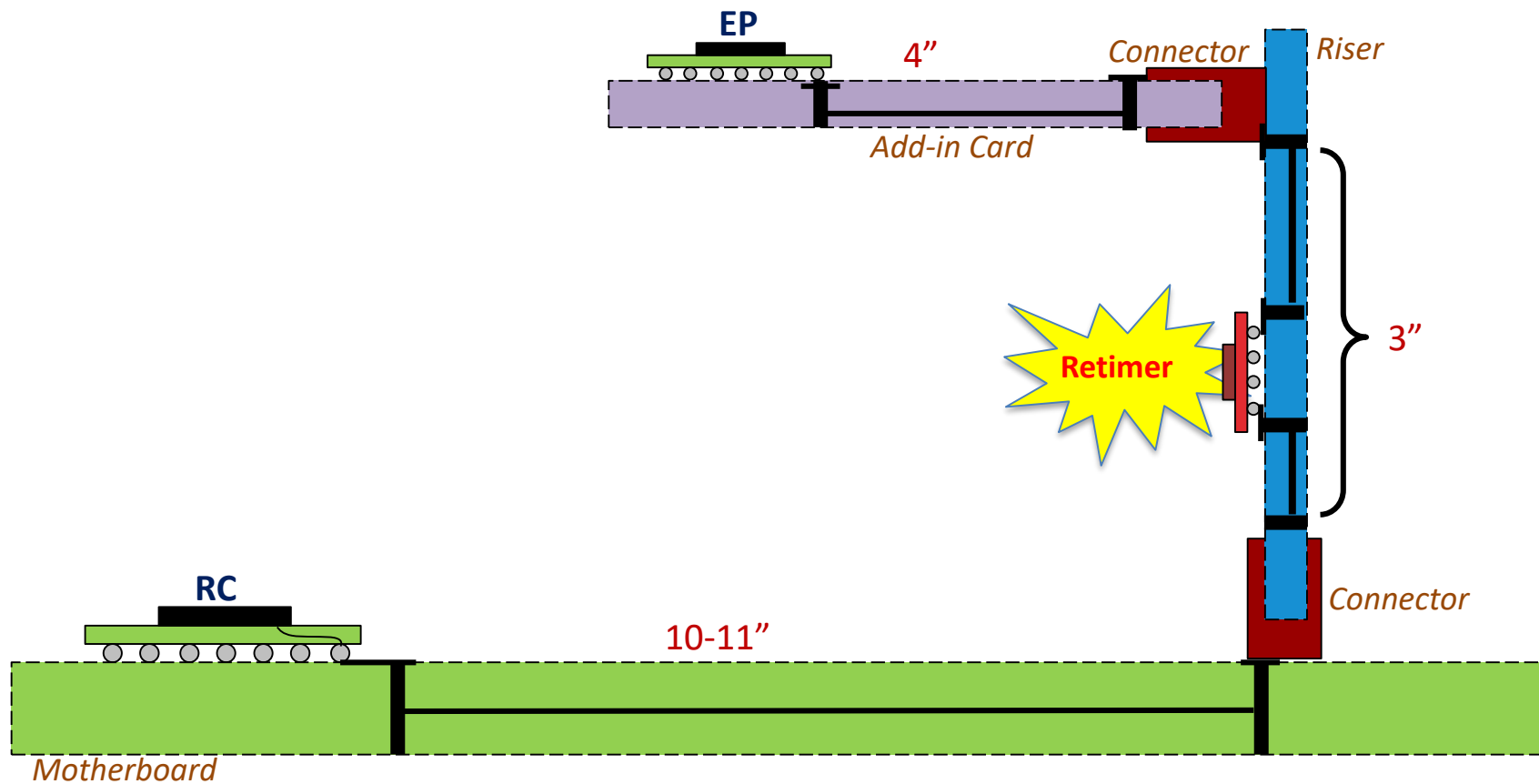
- **PCIe 4.0 is serious challenge from signal integrity point of view, especially for servers, storage and accelerator clusters platforms, where PCIe channel spending over relatively long traces, number of slots and connectors, cables and add-in cards**
- **Long and complicated channel exceeds -28db PCIe 4.0 spec defined limit**
- **Re-timers targeted to expend PCIe channel by implementing Ports with CDR and Equalization tuning capabilities**
- **While being transparent for Data/Transaction layers and above**

Range of Topologies

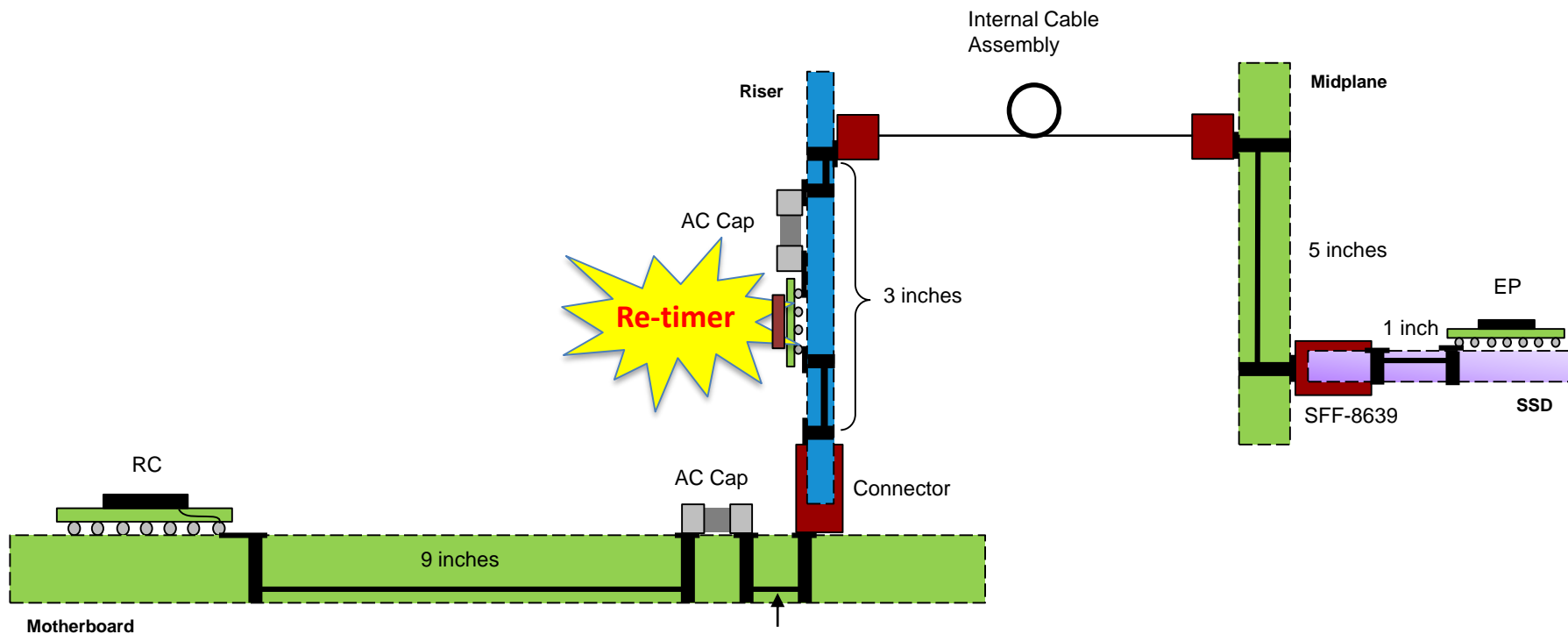


- **Server / Network FE / Storage Backplane / Accelerators Clusters / JBOF**
- **CPU ←RE-TIMER→ AIC**
- **CPU ←RE-TIMER→ Riser Card→ AIC**
- **CPU ←RE-TIMER→ Cable → Switch →AIC**
- **CPU → Switch ←RE-TIMER→ Cable →AIC**

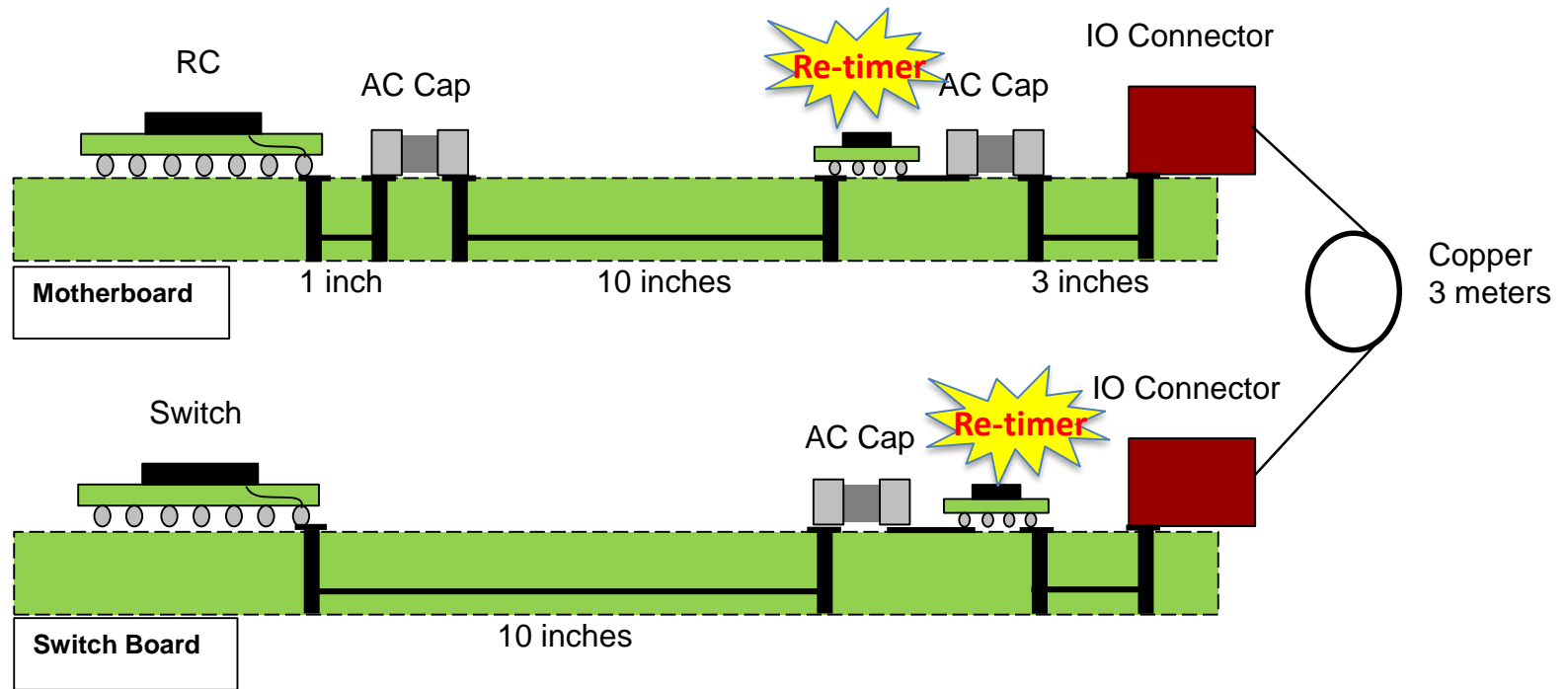
Re-timer Topology (1)



Re-timer Topology (2)



Re-timer Topology (3)



Re-timer Capabilities



- ✓ Physical Layer implemented
 - ✓ Phase 2/3 EQ procedure
 - ✓ EI and Inferred EI
 - ✓ RX impedance control on both side
- ✓ Totally transparent to DL and TL Layers
- ✓ TX: 3 stage FIR, full coefficient matrix
- ✓ Rx CDR
- ✓ RX DFE/CTLE
- ✓ Clock compensation by SKPOS add/remove
- ✓ Lane De-skew
- ✓ Lane numbering detection
- ✓ Lane polarity detection
- ✗ No standard LTSSM
- ✗ No CFG space
- ✗ NO DL/TL /AL
- ✗ No Error Handling . Errors are forwarded , including PHY /Symbol
- x4 /x8 /x16 **VVV**
- SRIS support on any pseudo port **V**
- Receiver Lane margining **V**
- Sideband registers rd/wr **V**
- Slave Loopback **V**

Agenda



- Introduction
- **System Integration Challenges**
- Re-timer Affect on Latency and Performance
- Conclusions

System Challenges



LATENCY

- Spec defines 64+ ns re-timer latency
- Up to two re-timers on path
- 250+ns additional latency on IO—CPU RTT
- *Latency and performance effect on next slides*

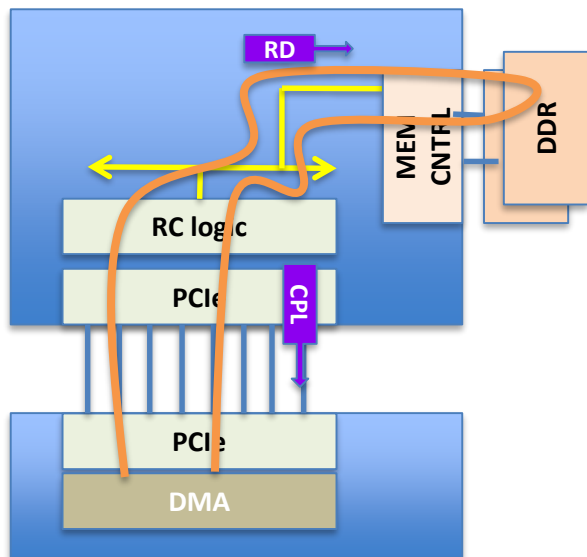
POWER

- Server CPU – 150W
- SSD – 30W
- x16 re-timer == 32 lanes → ~7W
- Additional 6-10% to platform power
- Heatsinks?

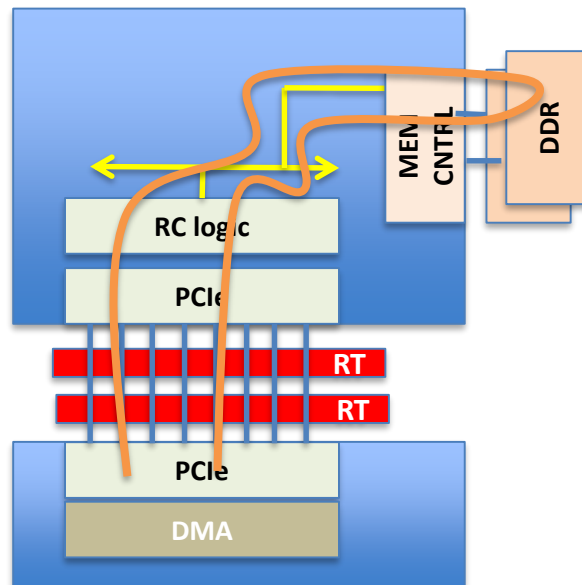
FOOTPRINT

- Additional component on board
- x2 link to SSD
- High density storage – M.2/BGA SSD

Latency

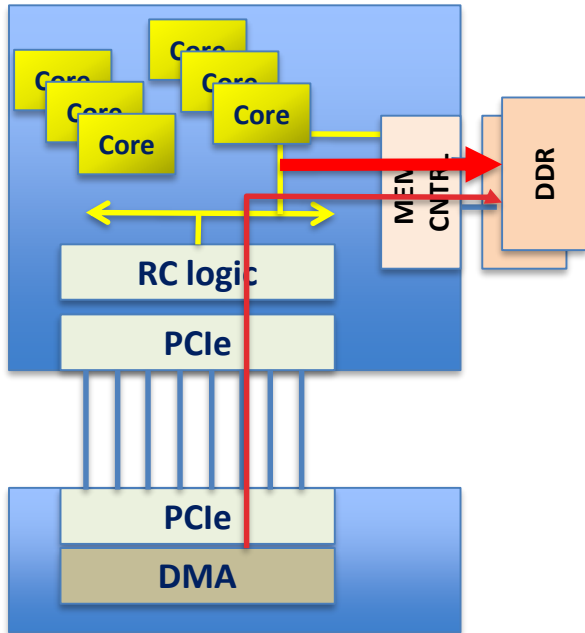


- **~250ns RTT**



- **~250ns RTT + 256ns two re-timers**
- **+100% to RTT**

Latency - Real Life Analyzes



RTT depend on dynamically change CPU load

RTT may vary from 250ns to **3usec**

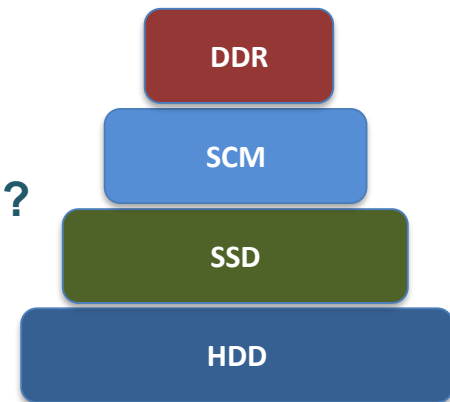
Generally IO devices optimize for latency <1usec

Performance may be degraded for longer RTT

PCIe SSD – latency of 100us → 20us ?

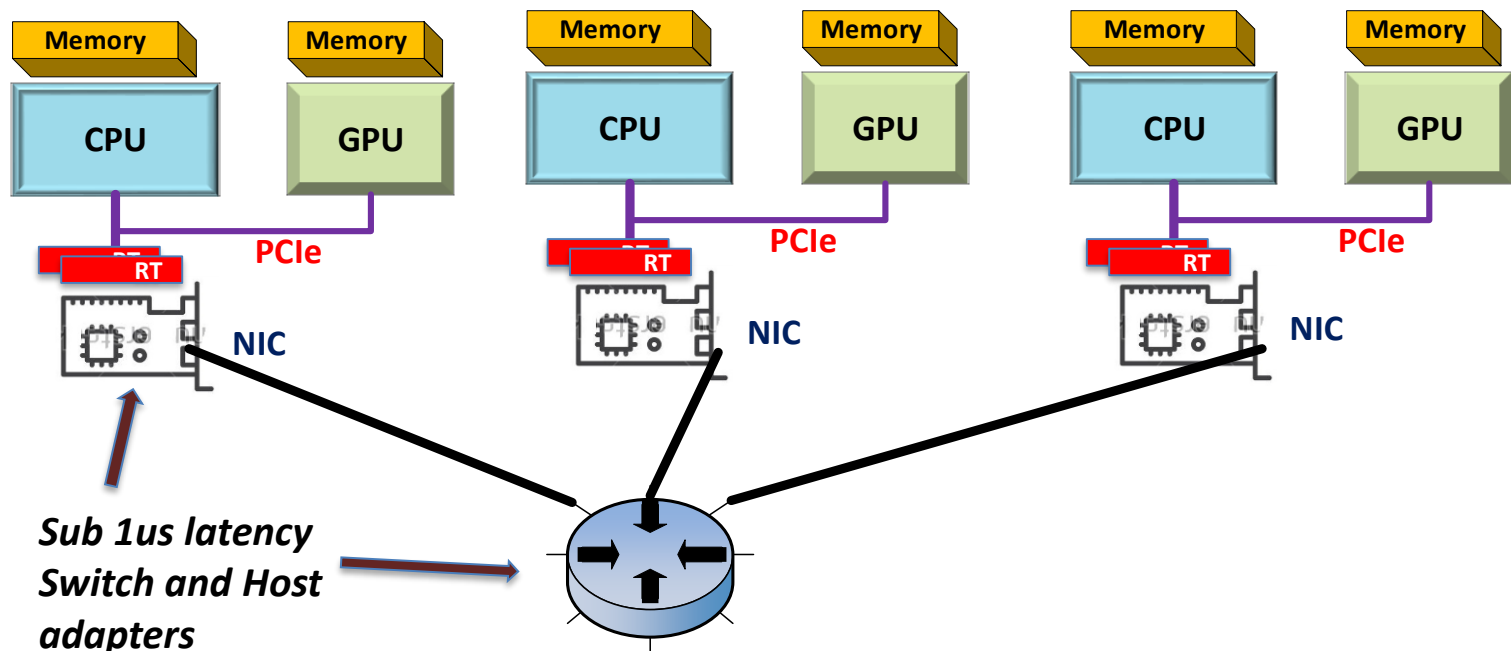
PCIe SCM ? Latency of 10us

PCIe FC HBA /NIC → HDDs (1-10ms)

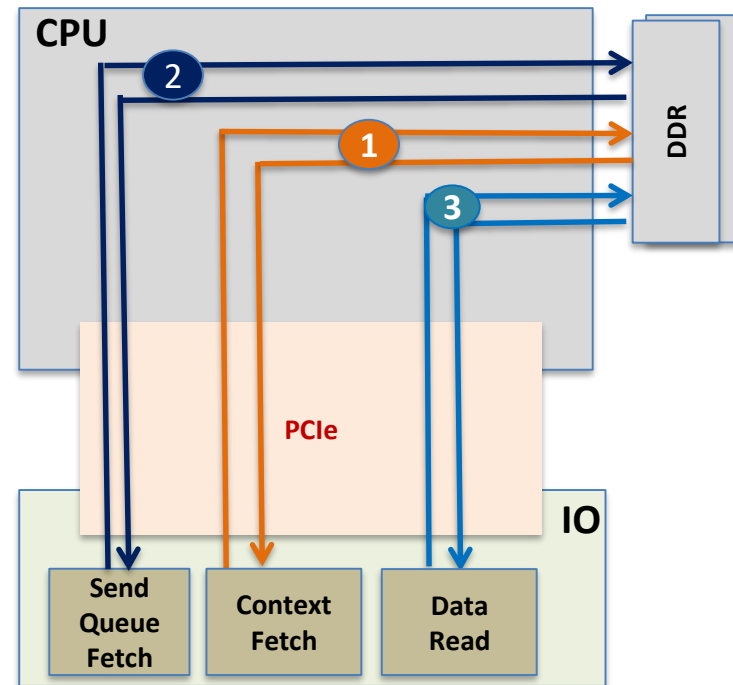
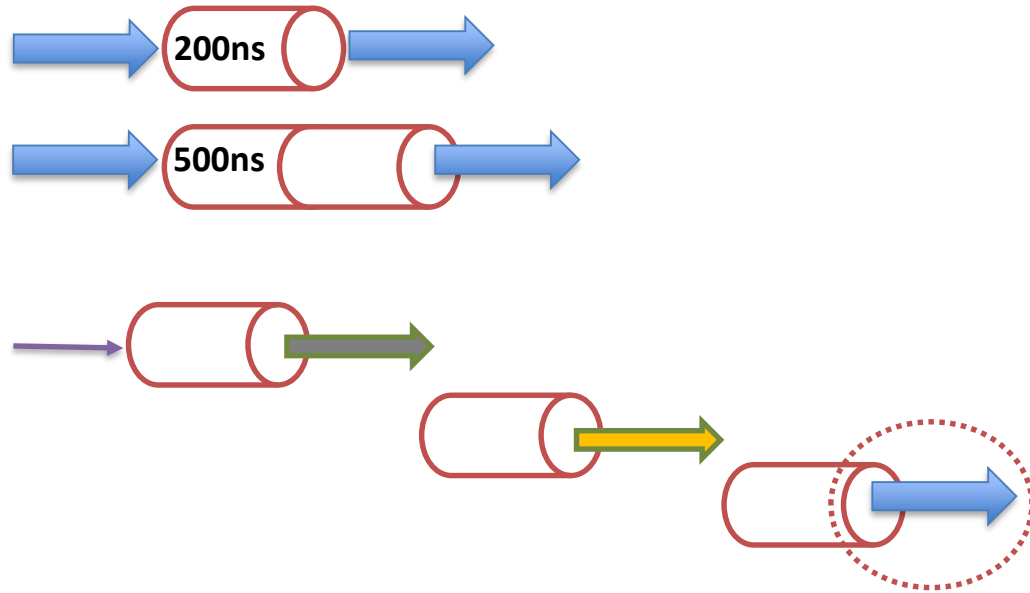


Latency in HPC Scenario

- Low latency networking (IB RDMA/RoCE) provides sub 1us latency
- Retimer latency will be increase this almost bare metal latency



Multi-stage IO Access



- Modern IO devices implementing Queues / Contexts / Tables in main memory
- Few accesses to main memory required for single chunk of data RX or TX operation
- Round trip latency may affect performance (IO/Packet rate)
- **IO try to prefetch Queues / Contexts / Data to mitigate RTT latency affect**

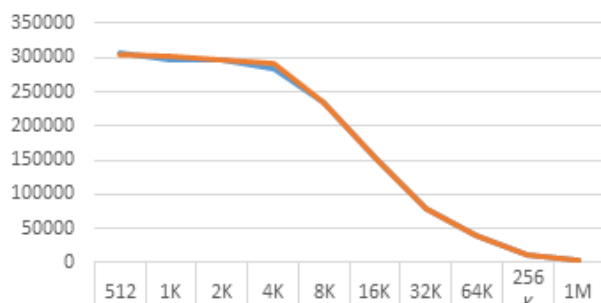
Performance vs. Latency (HCA)



- ~170ns latency PCIe used to emulate re-timers latency

Performance vs. Latency (SSD)

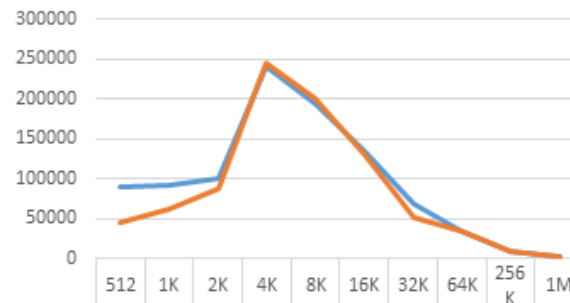
IOPS: Read Comparison



Direct connection Read	308222	295442	296517	283308	233885	153487	87434	40454	10336	2588
Connection SW Read	304233	300802	297812	290312	235546	153537	87964	40460	10348	2588

Direct connection Read Connection SW Read

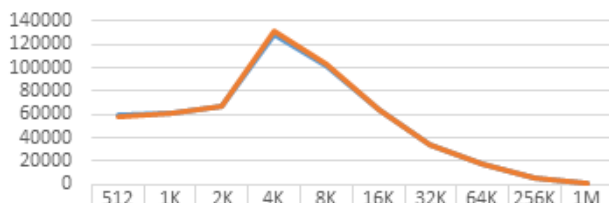
IOPS: Write Comparison



Direct connection Write	89179	90965	99622	242009	194452	134286	68533	33622	8680	2172
Connection SW write	45847	53002	86998	245927	200278	131005	51818	4114	8620	2168

Direct connection Write Connection SW write

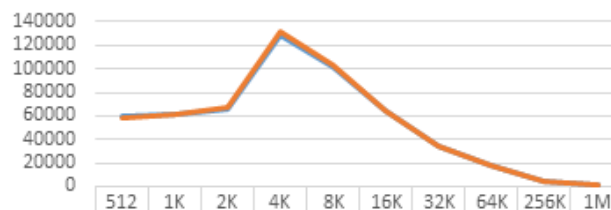
IOPS: Read&write mix-Write Comparison



Direct connection Read&write mix-Write	59189	60944	66753	129279	90176	64360	34147	17774	4672	1180
Connection SW Read&write mix-Write	58654	60796	66892	131329	90271	64322	33987	17823	4684	1176

Direct connection Read&write mix-Write
Connection SW Read&write mix-Write

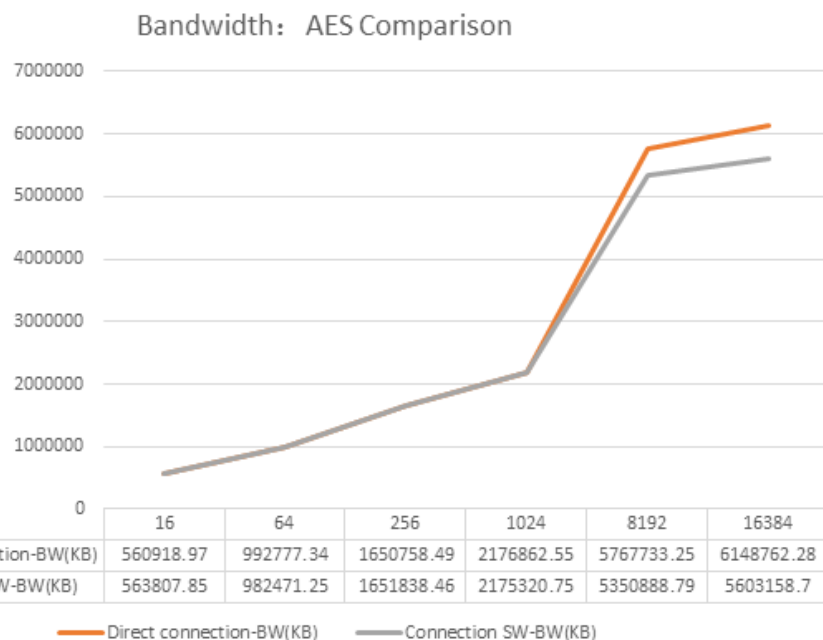
IOPS: Read&write mix-Read Comparison



Direct connection Read&write mix-Read	59201	60936	66590	129290	101745	64363	34175	17799	4680	1172
Connection SW Read&write mix-Read	58676	60797	66899	131368	102700	64335	34010	17850	4692	1168

Direct connection Read&write mix-Read
Connection SW Read&write mix-Read

Performance vs. Latency (ACCLs)



rsa2048	TYPE	sign	verify
	Direct connection-OPS	15422.6	182011.8
	Connection SW-OPS	15433.4	168018.1

- **Standard re-timer may add 64+ns latency in one direction. This latency is critical for low latency applications.**
- **Ultra low (sub 10ns) latency solution required.**
- **Data and Transaction layers are ready for affectional latency (Scaled FC/10b Tag).**
- **Fine tuning of IO PCIe side parameters may be needed.**
- **IO buffers increase may be requested to mitigate additional latency effect.**

**Thank you for attending the
PCI-SIG Developers Conference 2018.**

For more information, please go to www.pcisig.com

Don't forget to submit your feedback via the mobile app!

Download the **Crowd Compass** app and then search for **PCI-SIG Developers Conference** or entering the following URL into your mobile browser: <https://crowd.cc/s/1rKy0>

Enter event code: **DevCon2018**

Alternatively, access here: <https://crowd.cc/pcisig2018>

Note: Create an account within the app so Admin knows who to contact if selected as the prize winner.

**Each session feedback is provided is equivalent to 1 raffle entry (up to 11 sessions).
General survey feedback = 1 raffle entry.**

